# Optimizing Excited-State Electronic-Structure Codes for Intel Knights Landing



## Jack Deslippe

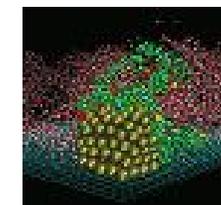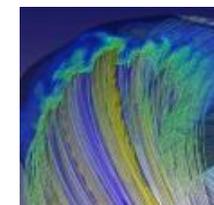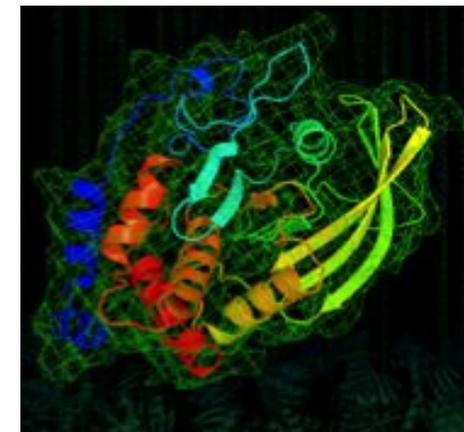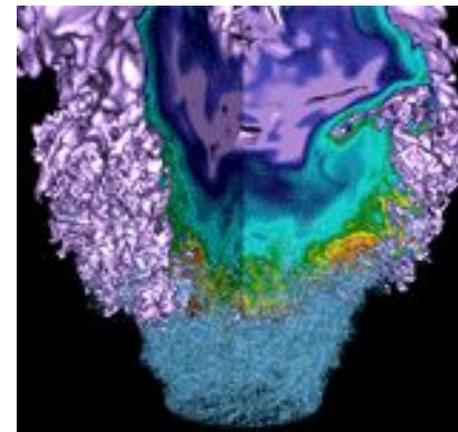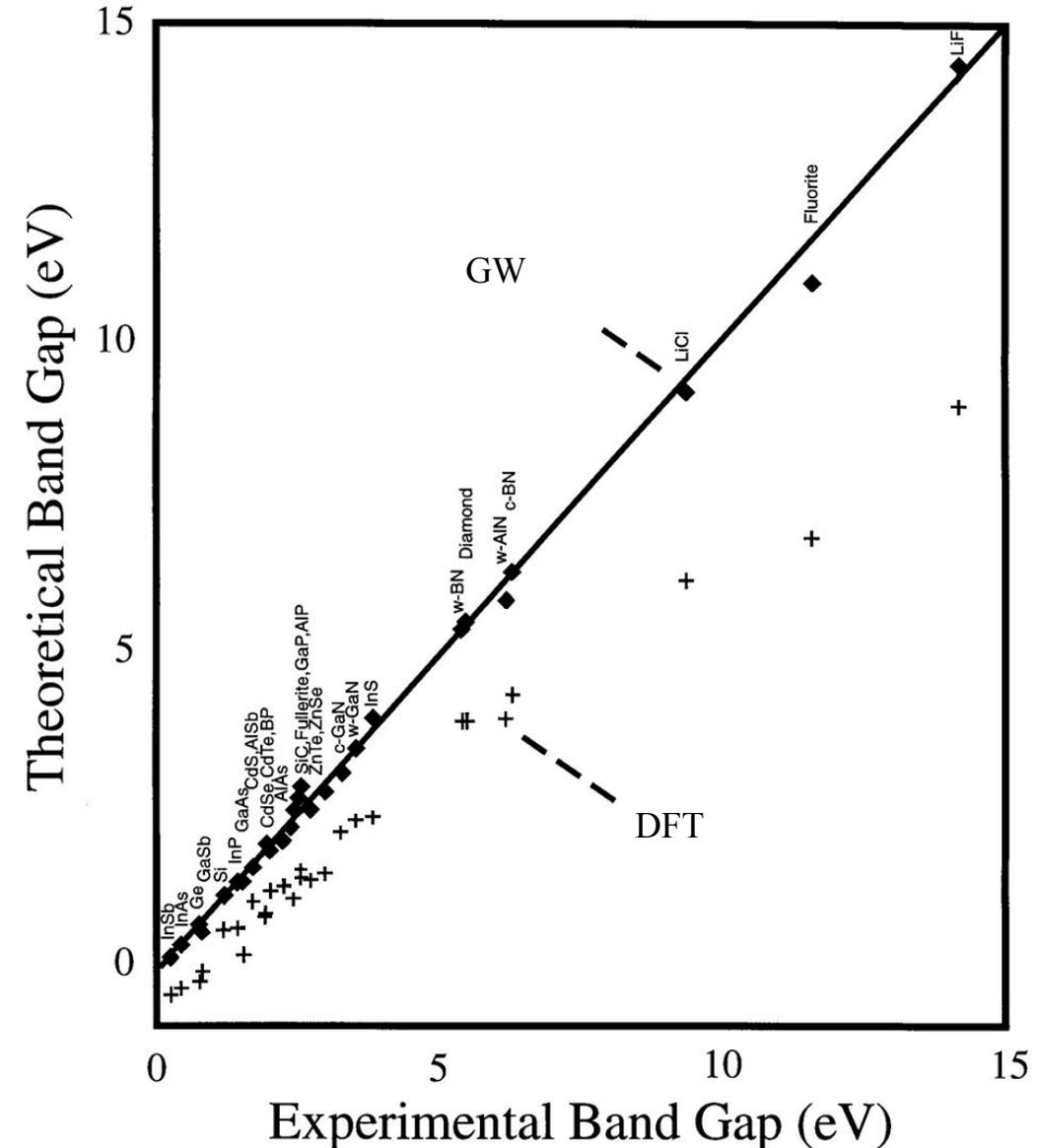**Application Performance Group**
**NERSC**

The "GW" method is an accurate approach for simulate the "excited state" properties of materials. Examples:

- What happens when you add or remove an electron from a system
- How do electrons behave when you apply a voltage
- How does the system respond to light or x-rays

GW is complementary to the widely used density functionally theory methods (DFT) which predict ground state properties of materials - i.e. the properties of the system associated with all particles in the lowest energy configuration.
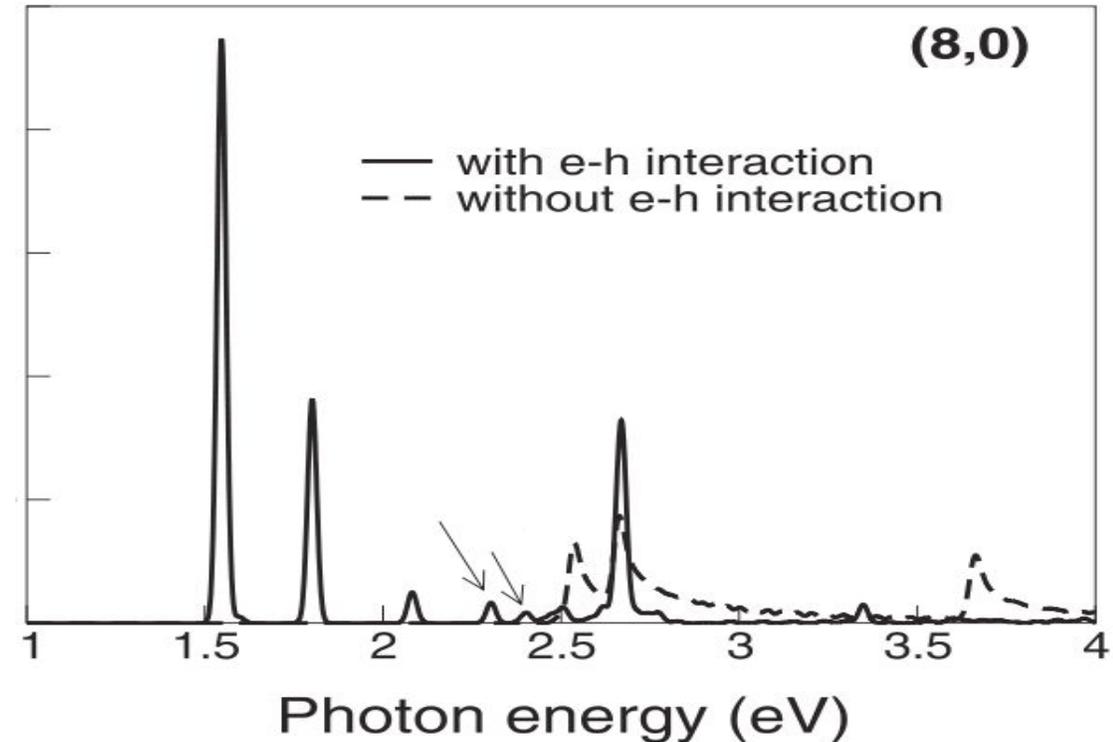
**Materials:**

InSb, InAs
Ge
GaSb
Si
InP
GaAs
CdS
AlSb, AlAs
CdSe, CdTe
BP
SiC
$C_{60}$
GaP
AlP
ZnTe, ZnSe
c-GaN, w-GaN
InS
w-BN, c-BN
diamond
w-AlN
LiCl
Fluorite
LiF

Many-body effects extremely important in **Excited-State properties** of Complex Materials.

Accurately describes properties important for:

- Photovoltaics
- LEDs
- Junctions / Interfaces
- Defect Energy Levels
- ….

**\*C.D. Spataru, S. Ismail-Beigi, L.X. Benedict, S.G. Louie. PRL 077402 (2004)**

**\*J. Deslippe, C.D. Spataru, D. Prendergast, S.G Louie. Nano Letters. 7 1626 (2007)**

Original code was MPI only. But, unlike DFT there are many layers over which you can exploit parallelism:

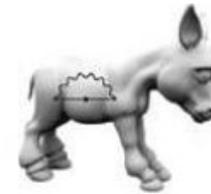$$\chi_{\mathbf{GG'}}(E) = \sum_{n}^{occ} \sum_{n'}^{emp} M_{nn'}^{*}(\mathbf{G}) M_{nn'}(\mathbf{G'}) \frac{1}{E_n - E_{n'} - E}$$

Band-pairs: $(n, n')$  Millions
Energies: $E$  Tens-Hundreds
Plane-Wave Basis Elements: $(\mathbf{G}, \mathbf{G'})$  - Millions

Much better suited exploiting levels of parallelism on HPC system like Cori. ~10,000 Nodes, 250 Threads per Node, 8 Wide Vectors.
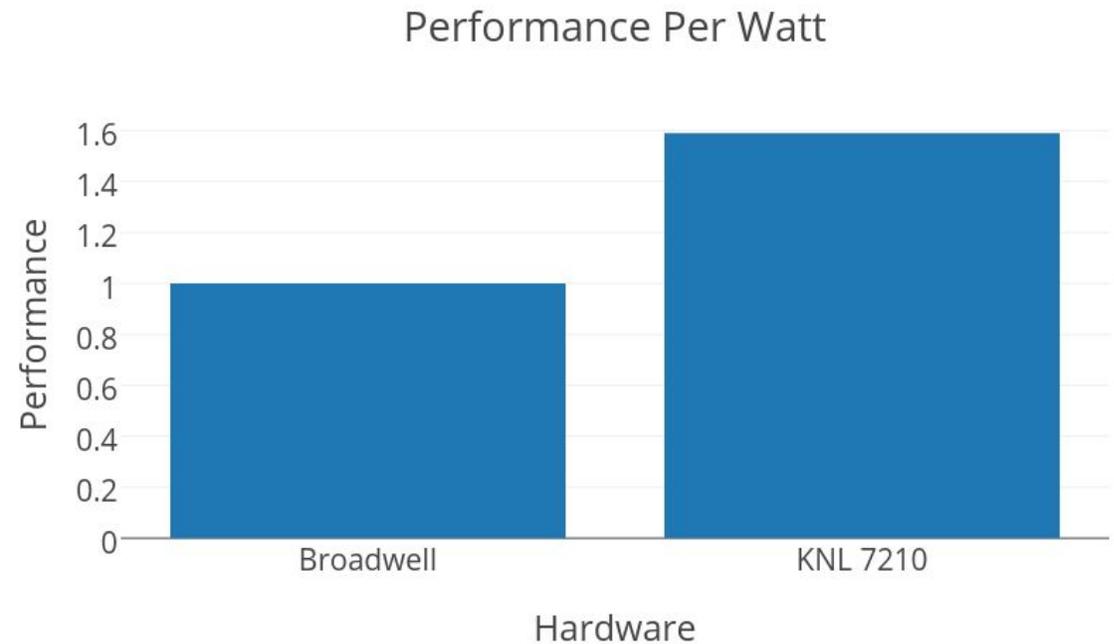
BerkeleyGW

# Don't Bury the Lead

**BerkeleyGW Running Well on KNL.**



Data From Sigma Benzene Runs on Single Node - Excluding IO.

# Computational Bottlenecks

**3 Main Computational Bottlenecks. Use roughly equal time for 500 atom systems.**

A. Compute transition probabilities (matrix-elements) for electrons from occupied to empty orbitals

B. Sum matrix-elements to form the overall material response function (polarizability)

C. Calculate the interacting electron energy from the polarizability

# Kernel A

Compute the transition probability between two electron states (orbitals):

$$M_{nn'}(\mathbf{G}) = \langle n| \, e^{i\mathbf{G}\cdot\mathbf{r}} \, |n'\rangle$$

Typically computed by FFT:

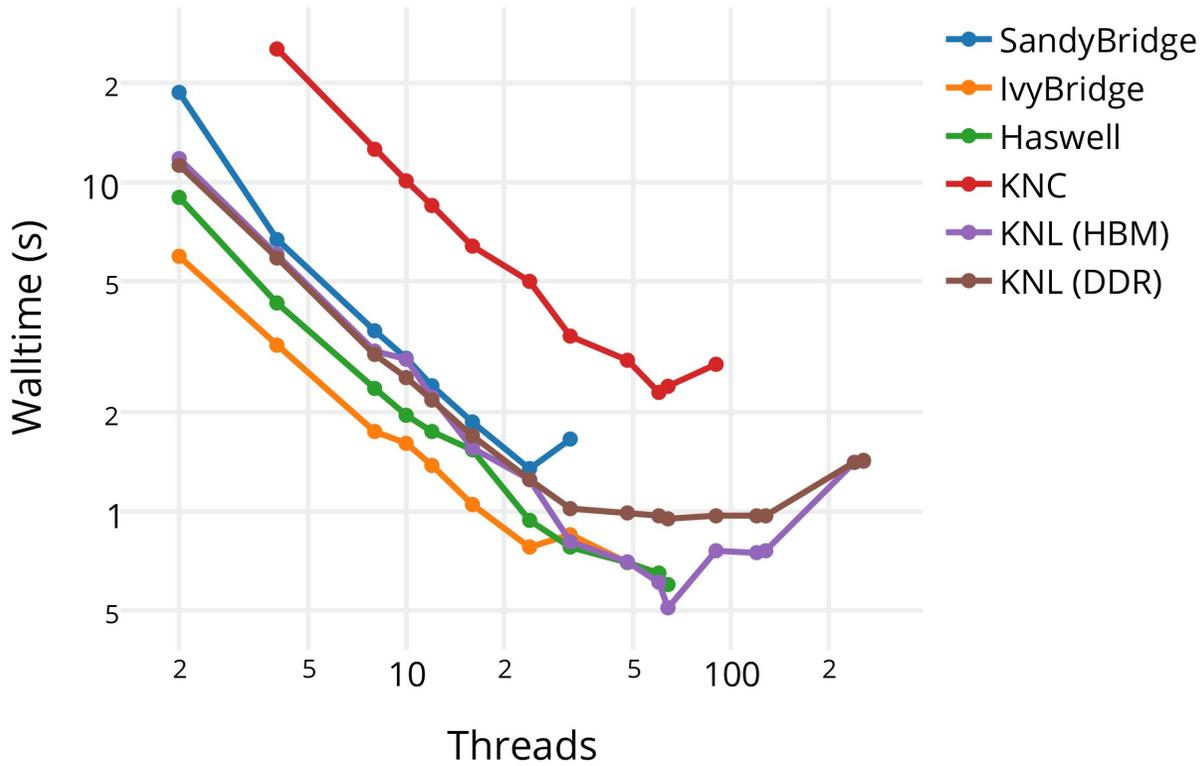$$M_{nn'}(\{\mathbf{G}\}) = \mathrm{FFT}^{-1}\left(\phi_n(\mathbf{r})\phi_{n'}^*(\mathbf{r})\right)$$

Must be done for all pairs of orbitals $n, n'$. Since the number of orbitals considered is proportional to number of atoms in calculation. The complexity of this step is O($N^3$logN).

We distribute ($n, n'$) via MPI and call threaded 3D FFT libraries (MKL) in the app. (Note, significantly more parallelism than local DFT)
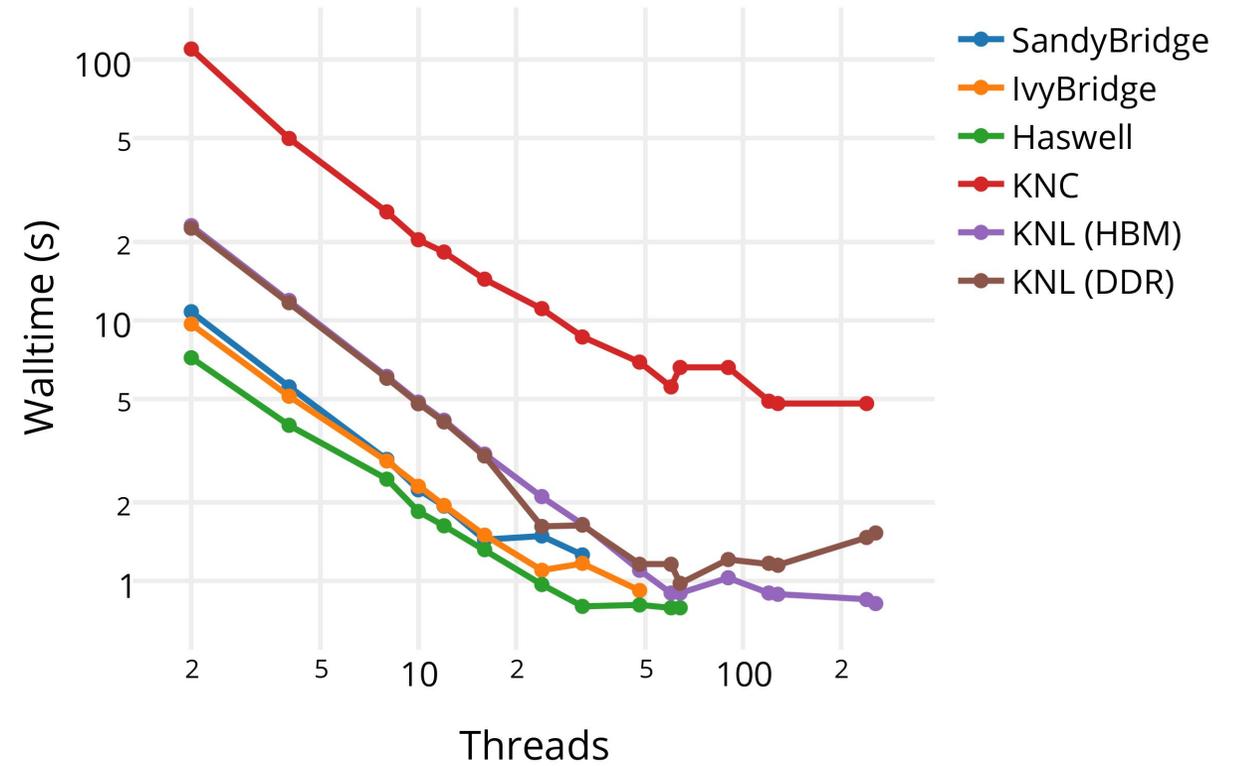
# Kernel A Performance (MKL)



Single FFT (960x960x480) Thread Scaling

Many FFT (135x135x135) Thread Scaling

# Related Improvements in Quantum ESPRESSO

When performing Hybrid Functional calculations within DFT, like within GW, you need to perform an FFT for each pair of orbitals.
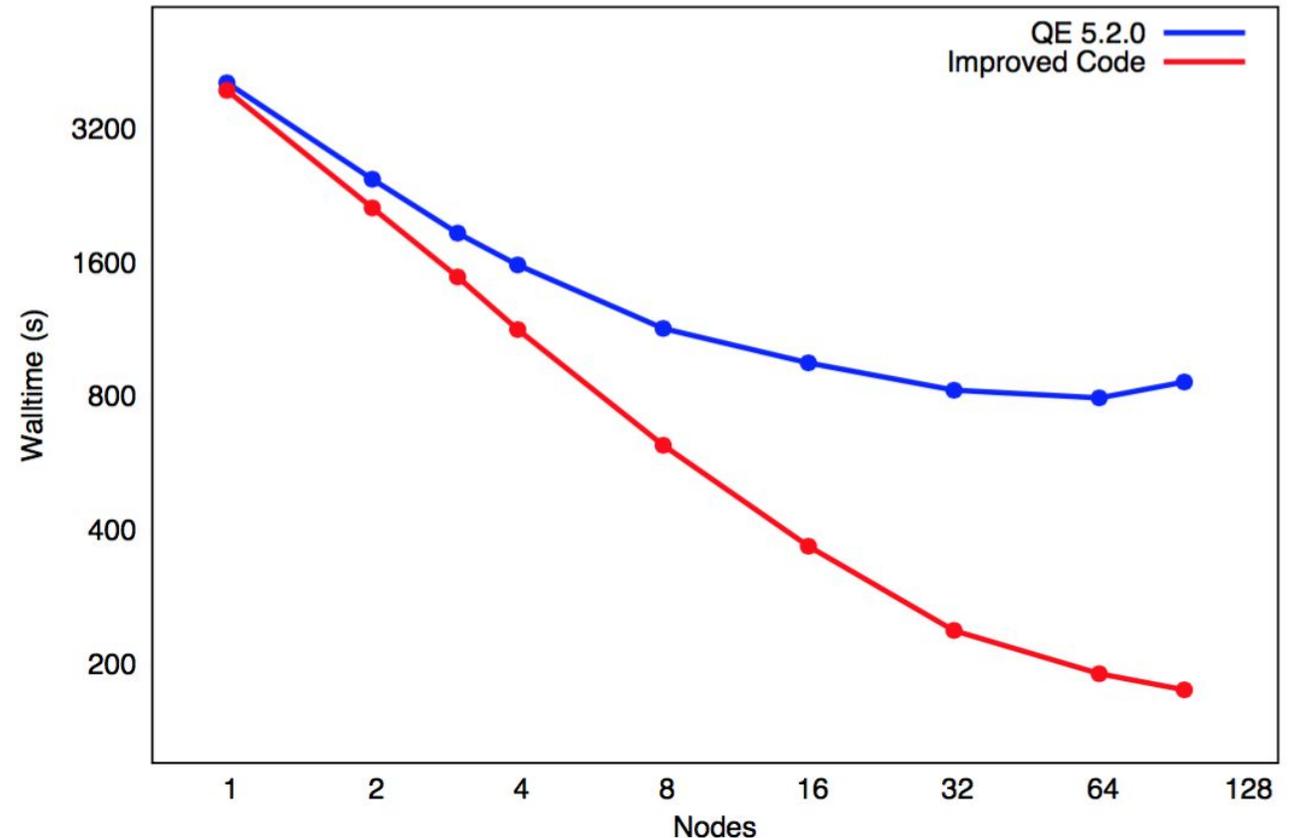
By default, code parallelizes each individual FFT.

We improve code to:

1. Parallelize over pairs of orbitals, before parallelizing individual FFTs

2. Allow simultaneous parallelism over orbitals for Hybrid calculation and other parameters for local calculation.

This leads to less communication, more work (complete 3D-FFT) on-node to parallelize.



Optimized QE on NERSC Edison

QE 5.2.0 (blue)
Improved Code (red)

We want to compute the electronic polarizability of the system:

$$\chi_{\mathbf{GG'}}(E) = \sum_{n}^{\text{occ}} \sum_{n'}^{\text{emp}} M_{nn'}^{*}(\mathbf{G}) M_{nn'}(\mathbf{G'}) \frac{1}{E_n - E_{n'} - E}$$

We can write this as a number of ZGEMM operations (one for each $E$):

$$\chi_{\mathbf{GG'}}(E) = \mathbf{M}^{*}(\mathbf{G}, (n, n'), E) \cdot \mathbf{M}^{\mathrm{T}}(\mathbf{G'}, (n, n'), E)$$

Where $\mathbf{M}$ is:

$$\mathbf{M}(\mathbf{G}, (n, n'), E) = M_{nn'}(\mathbf{G}) \cdot \frac{1}{\sqrt{E_n - E_{n'} - E}}$$

There are two steps. First, constructing $\mathbf{M}$, and second, performing the complex double-precision ZGEMM.

The complexity of this step is O($N^4$).
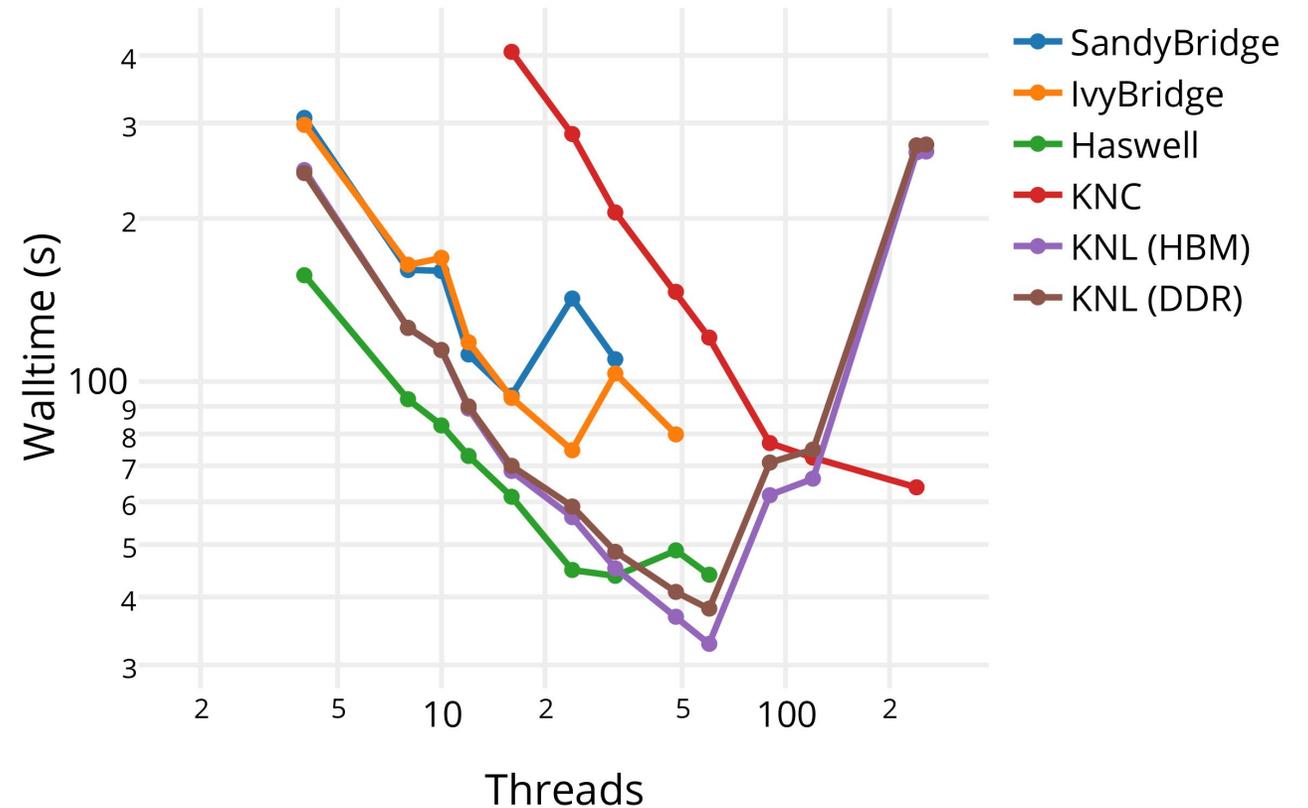
# Kernel B Performance

Little difference between MCDRAM and HBM performance - Only in the initialization/prep stage.

KNL overall performing 20% faster than Haswell.

No advantage of Hyper-Threading on Xeon or KNL seen.



Kernel B Thread Scaling

Compute the electronic energy as:

$$\Sigma_n = \sum_{n'} \sum_{\mathbf{GG'}} M_{n'n}^*(-\mathbf{G}) M_{n'n}(-\mathbf{G'}) \frac{\Omega_{\mathbf{GG'}}^2}{\tilde{\omega}_{\mathbf{GG'}} \left(E - E_{n'} - \tilde{\omega}_{\mathbf{GG'}}\right)} v(\mathbf{G'})$$
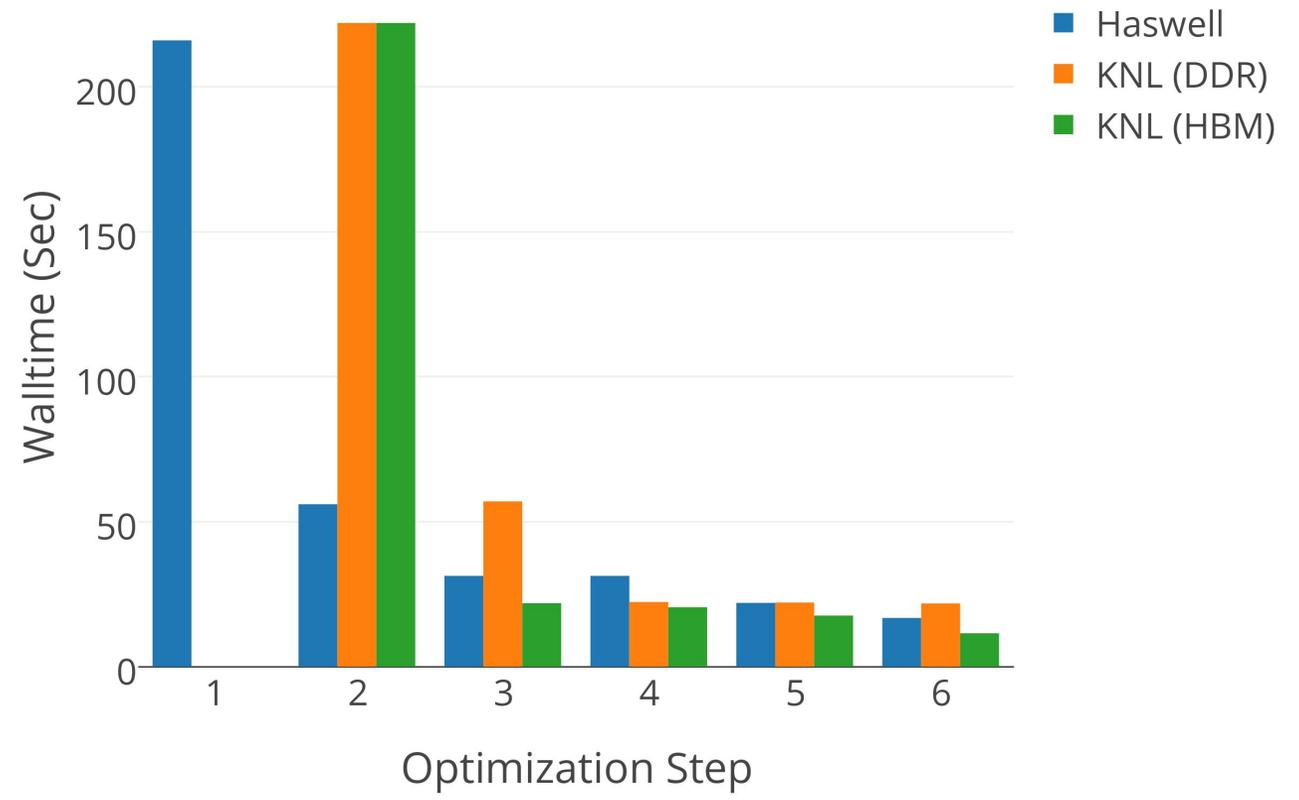
Here $\tilde{\omega}$ and $\Omega$ are complex double precision arrays derived from the polarizability. This is a tensor-contraction, matrix reduction type operation - performed by hand tuned code.

The complexity of this step for all $n$ is O($N^4$).

# Kernel C Optimization
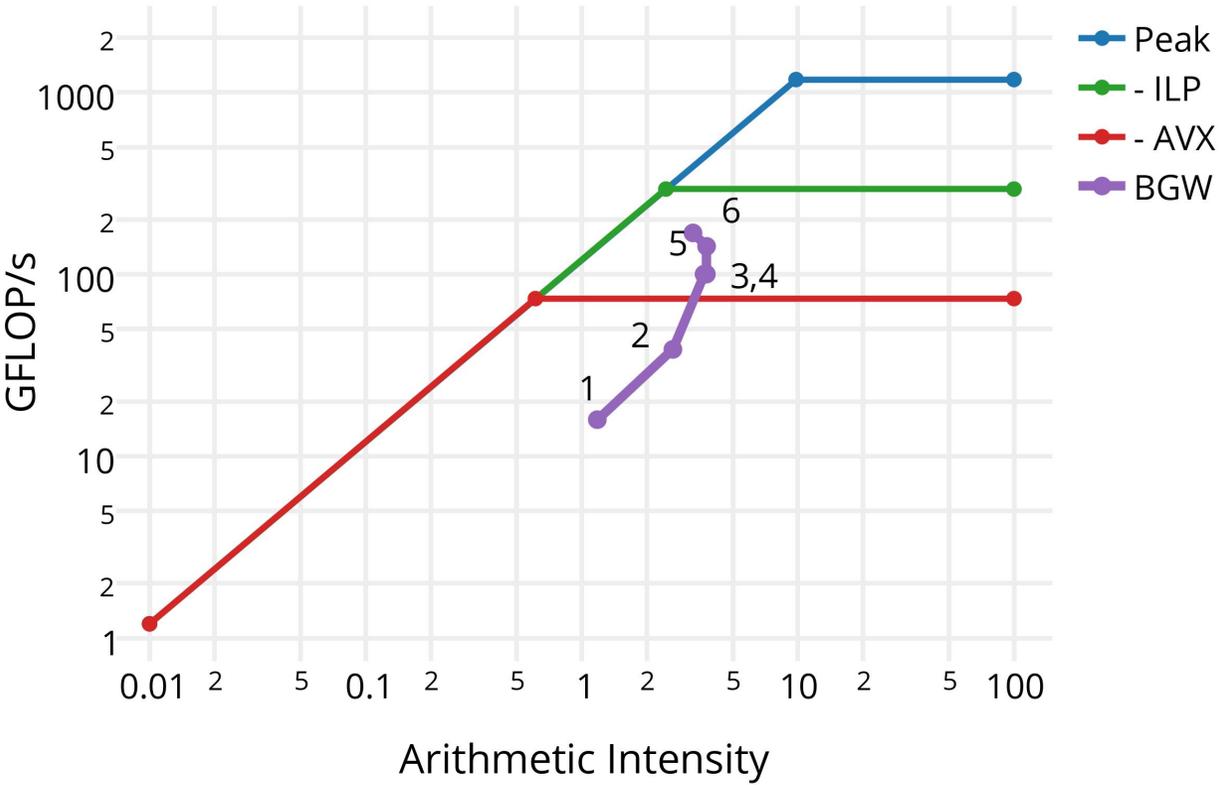
Optimization process for Kernel-C (Sigma code):

1. Refactor (3 Loops for MPI, OpenMP, Vectors)
2. Add OpenMP
3. Initial Vectorization (loop reordering, conditional removal)
4. Cache-Blocking
5. Improved Vectorization
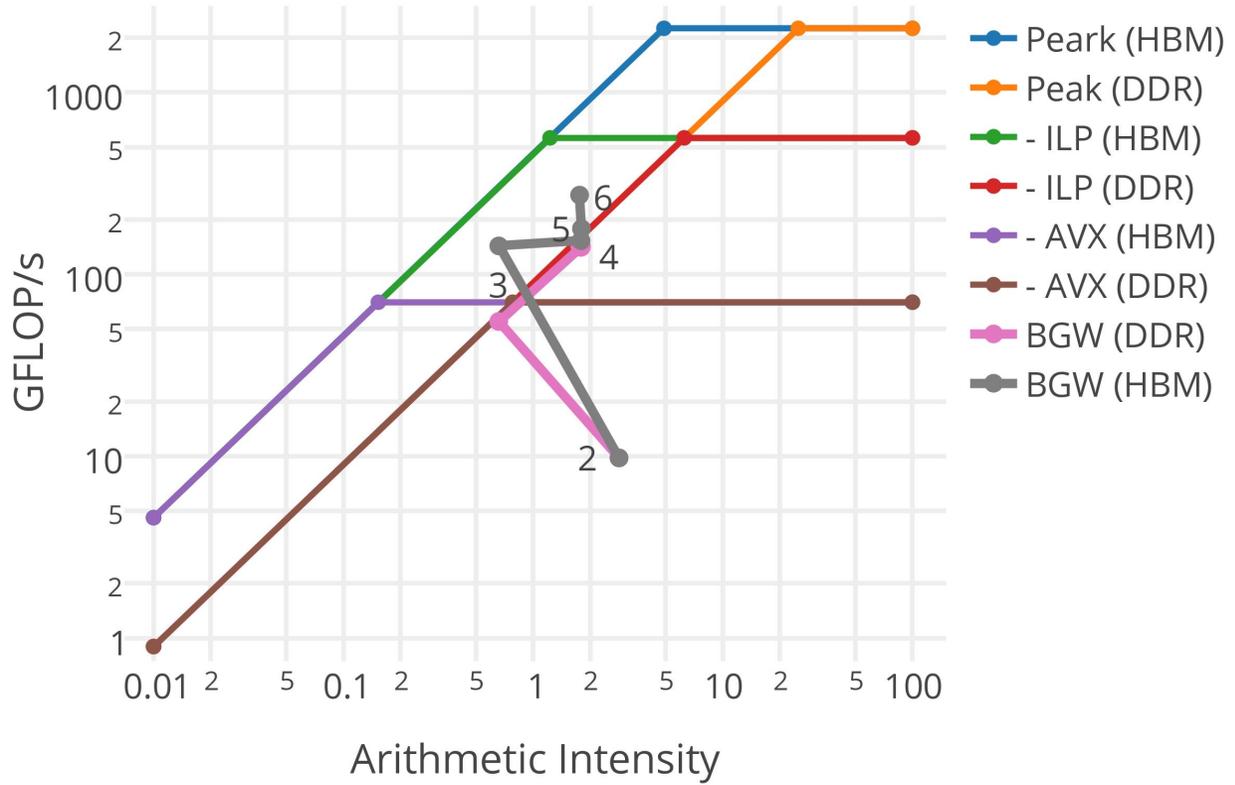6. Hyper-threading



Kernel C Optimization Process

Legend:
- Haswell
- KNL (DDR)
- KNL (HBM)

X-axis: Optimization Step
Y-axis: Walltime (Sec)

# Kernel C Optimization



Haswell Roofline Optimization Path

KNL Roofline Optimization Path

```
!$OMP DO reduction(+:achtemp)
  do my_igp = 1, ngpown
    ...
    do iw=1,3

       scht=0D0
       wxt = wx_array(iw)

       do ig = 1, ncouls

          !if (abs(wtilde_array(ig,my_igp) * eps(ig,my_igp)) .lt. TOL) cycle

          wdiff = wxt - wtilde_array(ig,my_igp)
          delw = wtilde_array(ig,my_igp) / wdiff
          ...
          scha(ig) = mygpvar1 * aqsntemp(ig) * delw * eps(ig,my_igp)
          scht = scht + scha(ig)

       enddo ! loop over g
       sch_array(iw) = sch_array(iw) + 0.5D0*scht

    enddo

    achtemp(:) = achtemp(:) + sch_array(:) * vcoul(my_igp)

  enddo
```

ngpown typically in 100's to 1000s. Good for many threads.

Original inner loop. Too small to vectorize!

ncouls typically in 1000s - 10,000s. Good for vectorization.

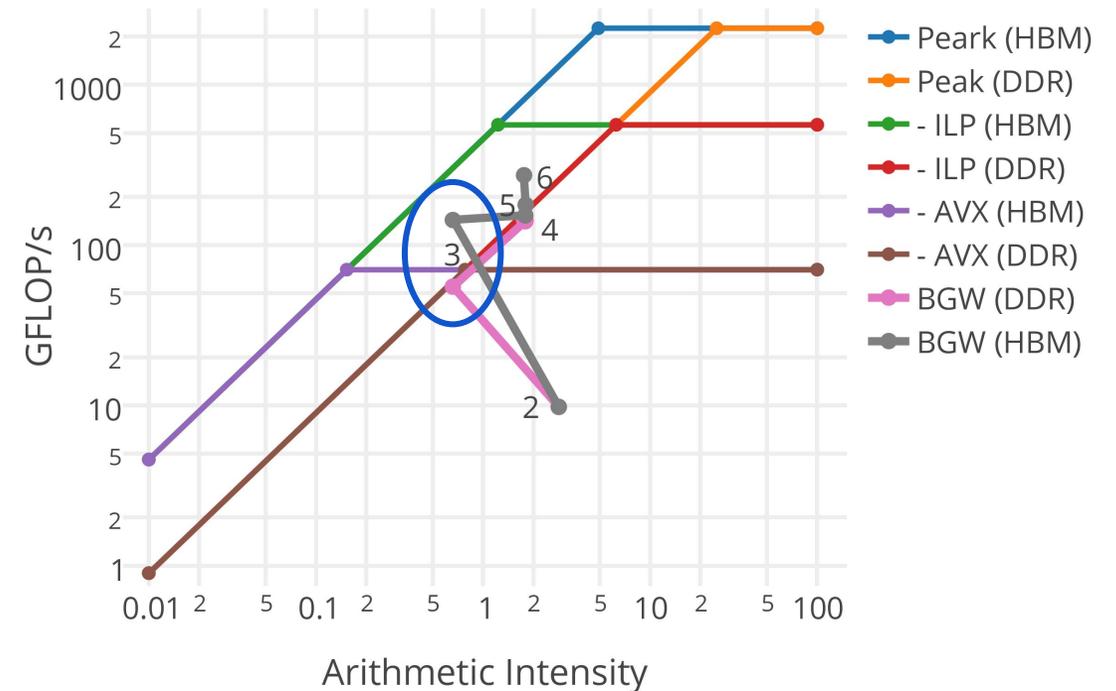Attempt to save work breaks vectorization and makes code slower.

Optimization process for Kernel-C (Sigma code):

1. Refactor (3 Loops for MPI, OpenMP, Vectors)
2. Add OpenMP
3. Initial Vectorization (loop reordering, conditional removal)
4. Cache-Blocking
5. Improved Vectorization
6. Hyper-threading

KNL Roofline Optimization Path



The loss of L3 on MIC makes locality more important.

```
!$OMP DO
do my_igp = 1, ngpown
    do iw = 1 , 3
        do ig = 1, igmax
            load wtilde_array(ig,my_igp)    819 MB, 512KB per row
            load aqsntemp(ig,n1)            256 MB, 512KB per row
            load I_eps_array(ig,my_igp)     819 MB, 512KB per row
            do work (including divide)
```

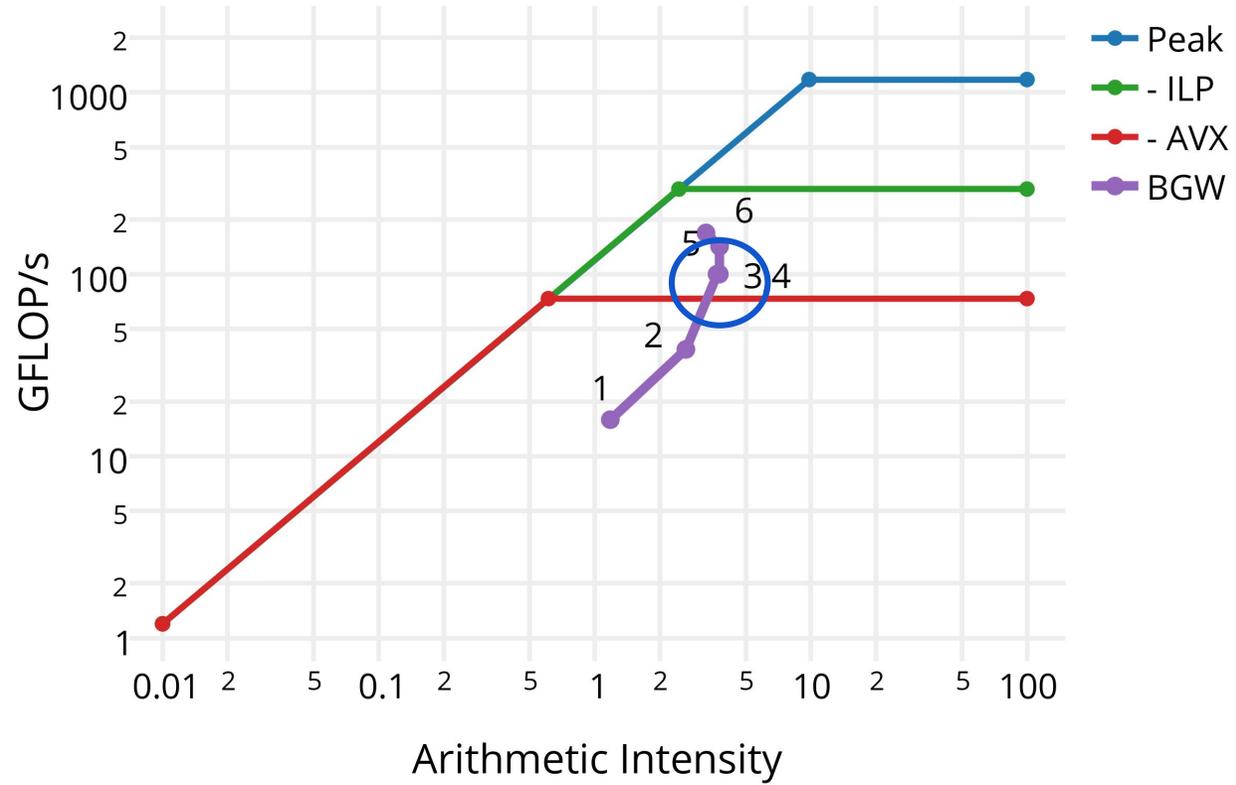Required Cache size to reuse 3 times:

1536 KB

L2 on KNL is 512 KB per core
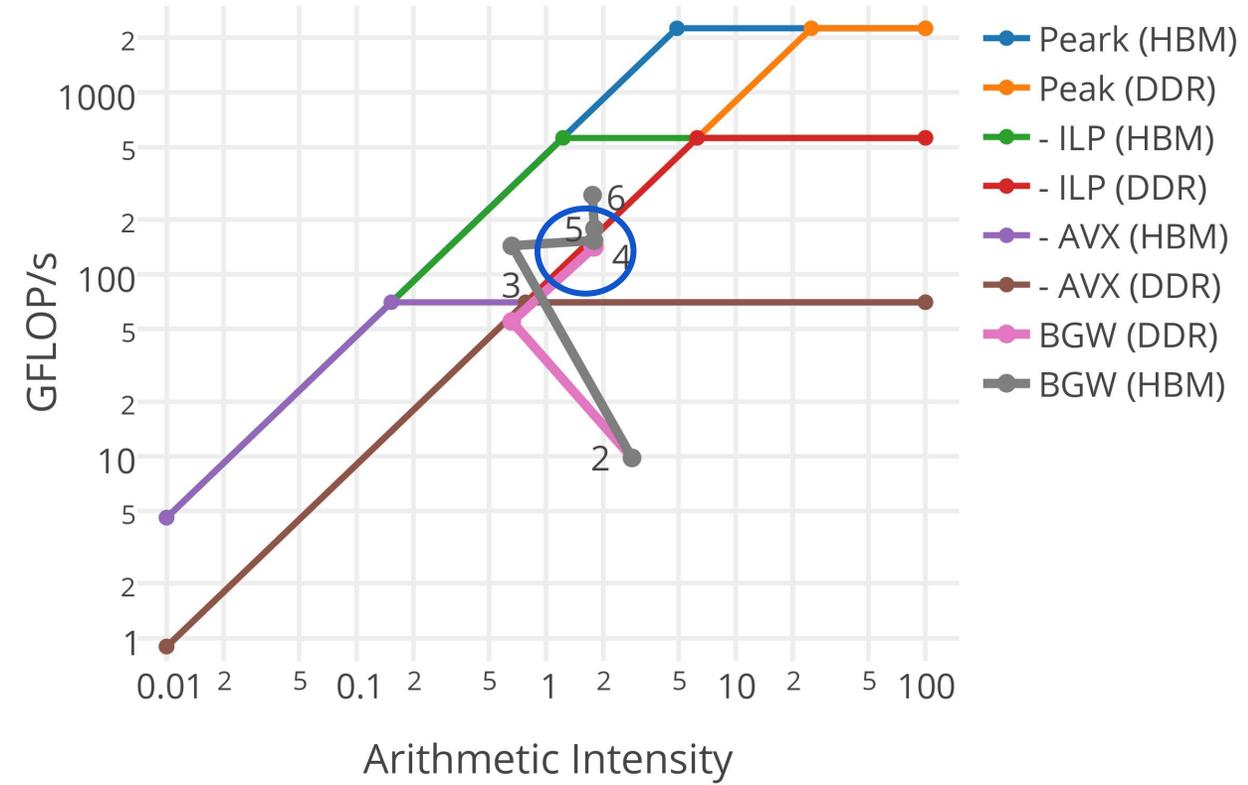L2 on Has. is 256 KB per core

L3 on Has. is 3800 KB per core

**Without blocking we spill out of L2 on KNC and Haswell. But, Haswell has L3 to catch us.**

```
!$OMP DO
do my_igp = 1, ngpown
   do igbeg = 1, igmax, igblk
      do iw = 1 , 3
         do ig = igbeg, min(igbeg + igblk,igmax)
            load wtilde_array(ig,my_igp)   819 MB, 512KB per row
            load aqsntemp(ig,n1)           256 MB, 512KB per row
            load I_eps_array(ig,my_igp)    819 MB, 512KB per row
            do work (including divide)
```

Required Cache size to reuse 3 times:

1536 KB

L2 on KNL is 512 KB per core
L2 on Has. is 256 KB per core

L3 on Has. is 3800 KB per core

**Without blocking we spill out of L2 on KNC and Haswell. But, Haswell has L3 to catch us.**

# Kernel C Optimization



Haswell Roofline Optimization Path

KNL Roofline Optimization Path

# SDE + Vtune, Why Complex Divides so Slow?

Found ~ 2x Instruction reduction from AVX to AVX512

However, found significant x87 instructions from 1/complex_number
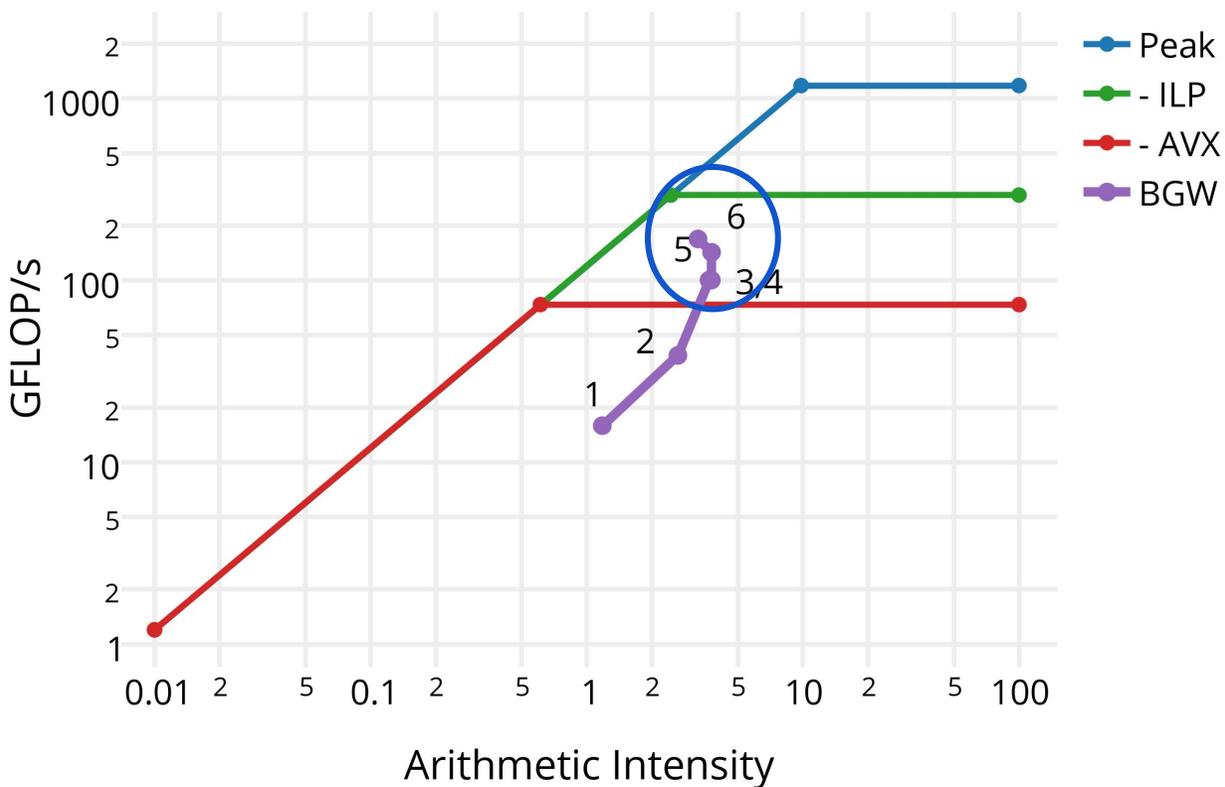


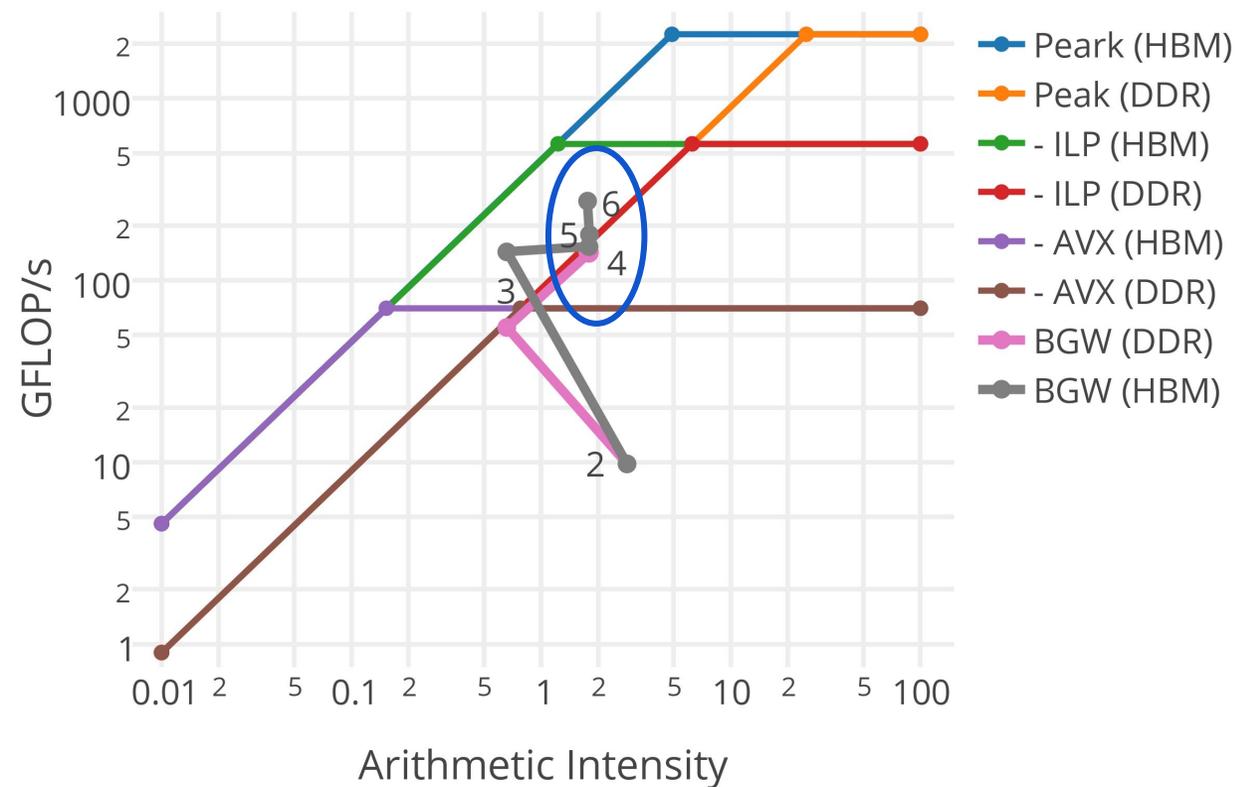Can significantly speed up by

a) Doing complex divide manually

Or

b) Using -fp-model fast=2

# Additional Speedups from Hyperthreading



Haswell Roofline Optimization Path
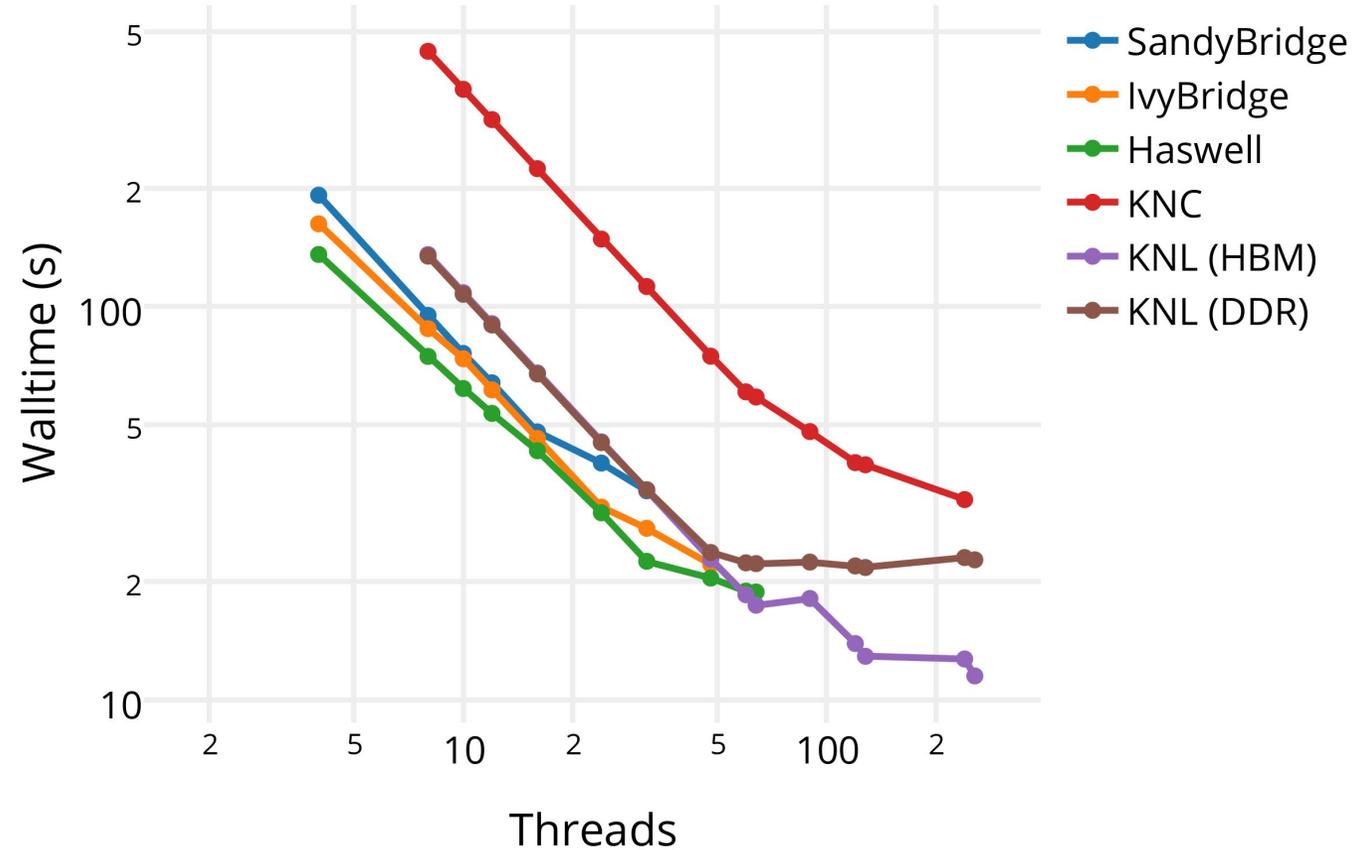
KNL Roofline Optimization Path

# Kernel C Thread Scaling

KNL DDR performance saturates at around 50 threads, becomes memory bandwidth limited.

KNL MCDRAM performance beats dual socket Haswell by 63%.
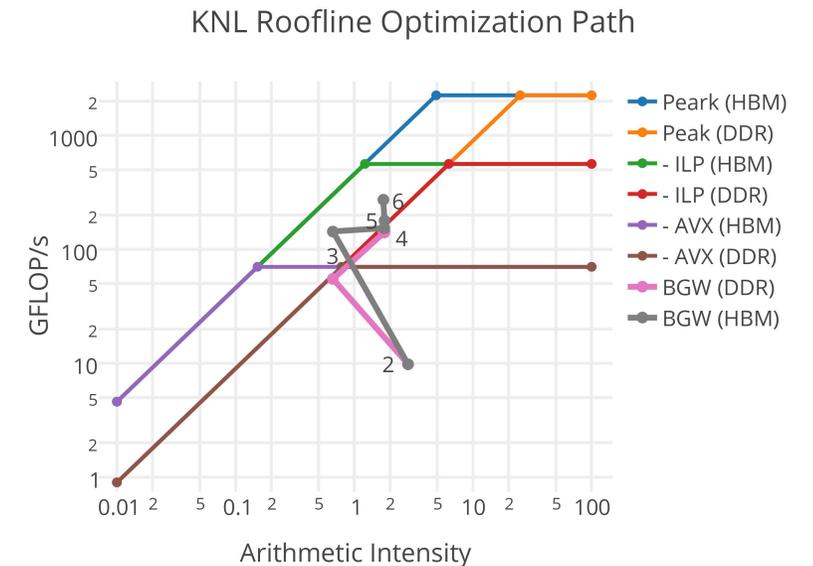


Kernel C Thread Scaling

Kernel A - FFTs show moderate speedups over dual-socket Haswell. Kernel B - ZGEMM and stream like operations show big speedups over Haswell

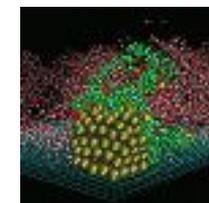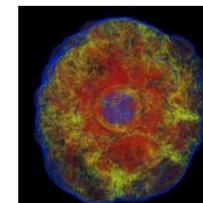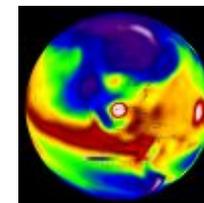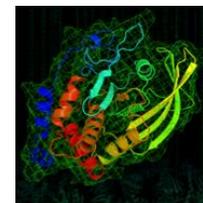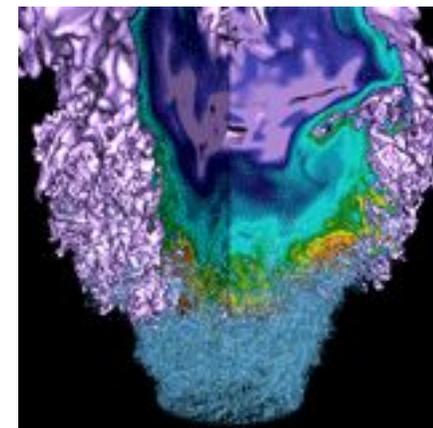Kernel C - Hand tuned matrix reduction operations show 60% speedup over haswell.

For algorithms with AI near roofline APEX (1-10), there is a rich optimization space that needs to be explored. Need all of:

- Thread scaling
- Vectorization
- Cache-Reuse
- Effective use of MCDRAM

Targeting Many-Core greatly helps performance back on Xeon.

KNL Roofline Optimization Path

# The End (Extra Slides)

$$[E_{n\mathbf{k}} - H_0(\mathbf{r}) - V_H(\mathbf{r})]\psi_{n\mathbf{k}}(\mathbf{r}) - \int \Sigma(\mathbf{r}, \mathbf{r}', E_{n,\mathbf{k}})\psi_{n\mathbf{k}}(\mathbf{r}')d\mathbf{r}' = 0$$

**The Good:**

Quantitatively accurate for quasiparticle properties in a wide variety of systems.

Accurately describes dielectric screening important in excited state properties.

**The Bad:**

Prohibitively slow for large systems.  Usually thought to cost orders of magnitude more time that DFT.

Memory intensive and scales badly.  Exhausted by storage of the dielectric matrix and wavefunctions.  Limited ~50 atoms.

$$\left[E_{n\mathbf{k}} - H_0(\mathbf{r}) - V_H(\mathbf{r})\right]\psi_{n\mathbf{k}}(\mathbf{r}) - \int \Sigma(\mathbf{r},\mathbf{r}',E_{n,\mathbf{k}})\psi_{n\mathbf{k}}(\mathbf{r}')d\mathbf{r}' = 0$$

**The Good:**

Quantitatively accurate for quasiparticle properties in a wide variety of systems.

Accurately describes dielectric screening important in excited state properties.

**The Bad:**

Prohibitively slow for large systems. Usually thought to cost orders of magnitude more time that DFT.

Memory intensive and scales badly. Exhausted by storage of the dielectric matrix and wavefunctions. Limited ~50 atoms.

# BerkeleyGW Towards Many-Core

★ In an MPI GW implementation, in practice, to avoid communication, data is duplicated and **each MPI task has a memory overhead.**

★ Users sometimes forced to use 1 of 24 available cores, in order to provide MPI tasks with enough memory. **90% of the computing capability is lost.**